# CRUK AND BIG DATA

**DR FIONA REDDINGTON**
**HEAD OF POPULATION RESEARCH FUNDING**

# Cancer Research UK

## WHAT WE SPEND

– We spent £341m on research in 2014/15

– The money we raise is spent on

- research
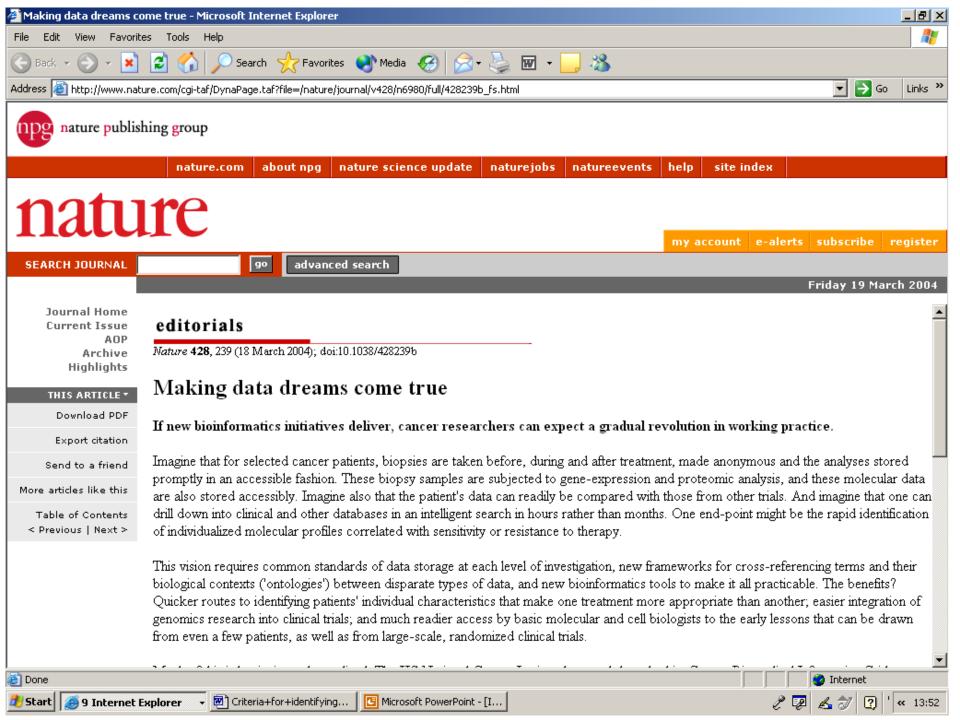- information
- advocacy and public policy

## WHERE WE WORK

– Cancer Research UK supports over 500 research groups

– We support research in about 40 towns and cities across the UK



CANCER
RESEARCH
UK

# Cancer Research UK's Ambition

Friday 19 March 2004

## editorials

# Making data dreams come true

**If new bioinformatics initiatives deliver, cancer researchers can expect a gradual revolution in working practice.**

Imagine that for selected cancer patients, biopsies are taken before, during and after treatment, made anonymous and the analyses stored promptly in an accessible fashion. These biopsy samples are subjected to gene-expression and proteomic analysis, and these molecular data are also stored accessibly. Imagine also that the patient's data can readily be compared with those from other trials. And imagine that one can drill down into clinical and other databases in an intelligent search in hours rather than months. One end-point might be the rapid identification of individualized molecular profiles correlated with sensitivity or resistance to therapy.

This vision requires common standards of data storage at each level of investigation, new frameworks for cross-referencing terms and their biological contexts ('ontologies') between disparate types of data, and new bioinformatics tools to make it all practicable. The benefits? Quicker routes to identifying patients' individual characteristics that make one treatment more appropriate than another; easier integration of genomics research into clinical trials; and much readier access by basic molecular and cell biologists to the early lessons that can be drawn from even a few patients, as well as from large-scale, randomized clinical trials.

# What do our researchers need to be able to do with data?

| Discover | Acquire | Store | Share | Mine |
|----------|---------|-------|-------|------|
|          |         |       |       |      |

- Collection of standardised datasets is becoming more commonplace for both research and clinical data

- Data sharing needs to improve

- Move towards provision of large scale datasets and infrastructure - not every research project needs to collect every data item

- Identify where re-use is feasible and where do we still need to fund collection of specific data items?

CANCER RESEARCH UK

# Generating Big Data is expensive

- PROSPECTIVE COHORT OF 500,000 PARTICIPANTS

- LIFESTYLE DATA, MEASUREMENTS AND BIOLOGICAL SAMPLES COLLECTED AT BASELINE

*COST = £90 MILLION FOR 5 YEARS*

- FOR £20 PER PARTICIPANT CAN COLLECT AND ANALYSE APPROX 50 BIOMARKERS – WILL TAKE 18 MONTHS

*COST = £10 MILLION TO GENERATE THE DATASET*

- CURRENT BID TO ADD IMAGING ENHANCEMENTS ON SUBSET OF 100,000 PEOPLE (MRI, DEXA, 3D ULTRASOUND)

*COST = £6 MILLION FOR PILOT, £26 MILLION FOR 5 YEAR STUDY*

- PLAN TO TRACK DISEASE OUTCOMES VIA LINKAGE TO ROUTINE COLLECTED MEDICAL DATA (E.G. HES, GPES, REGISTRY DATA)

*COST = ? TBC WITH DATA PROVIDERS*

THIS IS THE COST OF THE ESTABLISHMENT OF THE RESOURCE AND DATA GENERATION – IT DOES **NOT** INCLUDE THE COSTS OF ANY ASSOCIATED RESEARCH PERFORMED ON THE DATA

# Why is Big Data important in cancer research?

| Discovery | Clinical | Population |
|---|---|---|
| ICGC | YODA | UK Biobank |
| Human Genome project | CPRD | Cohort studies |
| Actionable Genome Consortium | CSDR.com | Routinely collected datasets (NHS and beyond) |
| | NPG Data Disclosure Project | |
| 100,000 genome, CRUK Strat Med | | |
| GA4GH | | |

And it's not just the Big Data.......

CANCER RESEARCH UK

# An example from clinical trials…

- Heterogeneity

    - Not all applications will be run through a CRUK Trials Unit (CTU)

    - 8 CRUK CTU's – all free to choose the data management software/systems they use

    - Multiple accreditation bodies (ECRIN, UKCRC, etc)

    - CTU's do not just work on cancer trials

    - CTU's do not just work on UK trials

Any guideline/policy/recommendation needs to be internationally compatible, system independent and non disease specific

# Future trials - What do we need?

- Data providers

  - Encourage applicants to think about entire data life cycle upfront when designing trial/study

    - Trial/cohort registration

    - Data discoverability – what data/what format?

    - Data accessibility – process for requesting data, how are decisions made, how many requests y/n, etc

  - Work to existing standards where they exist (e.g. CDISC, HL7, CONSORT, STROBE, etc)

  - Need common framework to work to where standards are not yet defined

CANCER
RESEARCH
UK

# Future trials - What do we need?

- Data requestors

  - What level of data is actually needed?

    - Raw data, CRF's, lay summaries

  - Is the data being requested for an appropriate purpose?

    - statistical analysis plan

    - valid scientific question

    - appropriate level of data requested

  - How will the results be made available and original team appropriately credited?

# What about trials that are closed/already underway?

- Need a rational way to prioritise:

  - What studies this should apply to

  - How far back should this apply?

  - How to resource – especially if funding for the trial has now ceased?

  - Does the data still exist!

# What about routinely collected datasets…

- Data have already been collected

- Data are discoverable

- Storage is taken care of by 3$^{rd}$ party provider

- So all our researchers need to do is access the data via agreed procedures/mechanisms

CANCER
RESEARCH
UK

**NHS**

Radiotherapy Dataset (RtDS) (collected by NATCANSAT - in principle a summary also goes to the registries)

Chemotherapy (Systemic Anti-Cancer Treatment Dataset) collected by the Chemo Intelligence Unit, sits in the Oxford Registry)

Cancer Outcomes and Services Dataset:
- Multidisciplinary team (MDT)
- Radiology
- Local imaging systems
- Patient administration system (PAS)

National Cancer Intelligence Network

LINKAGE

National Cancer Registration Service

Private healthcare: hospitals, hospices and histopathology

NHS Screening Programmes (bowel, cervix, breast)

Other UK cancer registries (when there is overlap)

GP records for deceased patients with no further

National PET-CT imaging contracts

Hospital Episode Statistics (HES - not cancer specific)

National Cancer Audit Data (bowel, head and neck, lung)

Cancer Waiting Times (collected by the Department of Health)

CRUK Stratified Medicine

Recurrent/Metastatic Breast

Death notification (cancer and non-cancer deaths of

Office of National Statistics:
- Calculate incidence and mortality
- De-duplicate data and provide basic quality checks

MDSCR (Minimum Dataset for Cancer Registration)

Clinical Practice Research Datalink

Health and Social Care Information Centre

GP Extraction Service (GPES - not cancer specific)

CPRD (anonymised, longitudinal medical records of patients from registered practices -

Diagnostic Imaging Data (not cancer

- There is a cancer registry for each **devolved nation**.
- In Scotland, patient data (excluding primary care data) is routinely linked by the Information Services Division. Information for researchers is provided by the Scottish Informatics Programme (SHIP).
- In Wales, the Health Information Research Unit operates the Secure Anonymised Information Linkage Databank (SAIL) which links and anonymises health data for research.
- In Northern Ireland, responsibility for health information sits with the Department of Health

# Informatics isn't just about data....

- We are entering an era where data generation is exploding

- Shortage of skills to undertake complex data linkage and analysis

- Methodological research needs development

- Understanding the science!

- Linking to data from other domains
  (e.g. National Pupil Database, Individual Learner Records, UCAS application records,
  Higher Education Statistics Data, Work and Pensions Longitudinal Study)

CANCER
RESEARCH
UK

# So what about preservation….

- Funder policies expect appropriate sharing, curation and preservation throughout data life cycle
  - responsibility of data custodians
  - does not always align with funding cycles

- Linking to data from other domains is important – other models to learn from?

- This requires partnership

# How can we partner effectively?

**THANK YOU!**

**cruk.org**

Dr Fiona Reddington

Head of Population, Prevention and Behavioural Research

020 3469 5324

**fiona.reddington@cancer.org.uk**

Angel Building, 407 St  John Street

London EC1V 4AD